

A31075-PCT USA 070050.1278

**IN THE UNITED STATES PATENT AND TRADEMARK OFFICE**

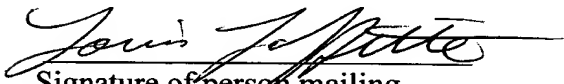
Applicant : CHANG, Shih-Fu et al.  
Serial No. : To be assigned  
Filed : November 4, 1997  
For : VIDEO SIGNAL FACE REGION DETECTION

**EXPRESS MAIL CERTIFICATION**

Express Mail Mailing No. EJ339571055US

Date of Deposit - May 1, 2000

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 C.F.R. §1.10 on the date indicated above and is addressed to: Box PCT, Assistant Commissioner for Patents, Washington, D.C., 20231.

  
Signature of person mailing  
correspondence

Name of person mailing correspondence: Louis Laffitte

DescriptionVideo Signal Face Region DetectionI. INTRODUCTION

This invention relates to video signals representing pictorial information, and particularly to a method for identifying portions of such video signals representing the face of a person.

In recent years, techniques have been proposed that allow users to search images by visual features, such as texture, color, shape, and sketch, besides traditional textual keywords. Algorithms have also been proposed to explore automatic detection of low level visual features directly from the compressed  
5 domain, based on the synergy between feature extraction and compression.

The human face is an important subject in image and video databases, because it is a unique feature of human beings, and is ubiquitous in photos, news video, and documentaries. The face can be used to index and search images and video, classify video scenes (e.g., anchorperson shots in news video), and segment  
10 human objects from the background. Therefore, research on face detection is critical in image and video database searching applications.

Although face detection is related to face recognition, the problem addressed is a little different from those in traditional face recognition applications. Prior work on face recognition has been focused on digital images taken in highly  
15 constrained environments. Strong assumptions are used to make the task more tractable. For example, there is usually just one front-view face in the center of the image, the head is upright, the background is clean, no occlusion of faces exists, no glasses are worn, and so on. The existence and locations of human faces in these images are known a priori, so there is little need to detect and locate faces. Face  
20 recognition has been an active research field for more than twenty years. Techniques for face recognition have been proposed using neural nets, deformable template matching, and Karhunen-Loeve expansion (i.e., eigenfaces). All these

002260 42505650

methods, as well as traditional ones, can be roughly classified into two broad categories, namely geometric-feature-based matching, and template matching.

In image and video databases, however, there is generally little or no constraint on the number, location, size, and orientation of human faces in the image or video scenes. The background of these images and video scenes is usually complex. Thus, face detection becomes important and challenging before the indexing, search, and recognition of the faces can be done.

Recently, work has begun on face detection in unconstrained images. The phrases face detection, face location, and face localization, are used interchangeably in literature. For consistency we use the phrase face detection herein. Govindaraju et al proposed a model-based approach where the face is defined as inter-connected arcs that represent chins and hairline. The arcs are extracted using low level computer vision algorithms, and are then grouped based on cost minimization to detect candidate face regions. Ten images are tested without any miss. However, false alarms are often generated. Burl et al presented a face detection system in which local feature detectors are coupled with a statistical model of spatial arrangement of facial features. A set of facial features (e.g., eyes, nostrils, nose-mouth junction) is detected first, and is used to form constellations of feature sets. A ranking technique with a statistical model of spatial arrangements of these features is then applied to detect possible face regions in the constellations of feature sets. To facilitate the detection process, an intelligent search scheme in the constellations is used. Evaluation on a database of 150 images (quasi-frontal, under the same lighting condition) indicates a correct detection rate of around 84%. A low complexity algorithm to detect and track face regions was proposed by Eleftheriadis et al. for model-assisted coding of low-bit-rate teleconferencing video. The algorithm is a three-step hierarchical procedure utilizing the fact that the human face outline is roughly elliptical. Yang and Huang proposed a hierarchical knowledge-based algorithm to detect human faces in a complex background. The algorithm consists of three levels. The higher two use mosaic images of different resolutions. The third one extracts edges of

facial components. Domain knowledge and rules are applied at each level. A detection rate of 83% is reported (50 faces from 60 512x512 images) with 28 false alarms. The run time of face detection is 60-120 seconds on a SUN Sparc 2 workstation. A neural network-based face detection system was reported  
5 comprehensively by Rowley et al. A set of neural network-based filters is first applied to an image at several scales. An arbitrator is then used to combine the filter outputs. The algorithm is able to detect 90.5% of the faces in 130 images from three different sources, many of which contain multiple faces. Computational complexity is high because the neural networks have to process many small local  
10 windows in the images. Wavelet transform domain has been explored for face detection as well. Venkatraman and Govindaraju used zero-crossings of a wavelet transform at different scales to extract local facial features. These features are then combined in a model matching stage to detect faces.

It is an object of the present invention to provide a new method for  
15 detecting face regions in a video signal.

It is a further object of the invention to provide a method which has a short processing time and provides a high probability of identifying face regions with low false alarm rates.

It is a further object to provide a method that can detect face regions in  
20 MPEG compressed video.

#### Summary of the Invention

In accordance with the invention there is provided a method for identifying face regions in a color image. The method includes providing image representative data, including data representative of chrominance in incremental portions of the  
25 image. The chrominance data for each image portion is compared to values known to be representative of skin tone to distinguish skin tone image portions from other image portions. The shape of regions having contiguous skin tone image portions is compared to at least one template consistent with the shape of a human face image to identify possible face regions.

In a preferred embodiment, the templates are rectangular and have vertical to horizontal aspect ratios between about 1 and 1.7. A further step may be provided of comparing the spatial frequency of data representing luminance to threshold values and eliminating false face regions having spatial frequency components below the threshold.

The method is advantageously used in connection with video frames in compressed MPEG signal format and using image portions corresponding to macroblocks of the I frames. Preferably the comparison of the chrominance signal is made using the DC component thereof. The shape comparison can be performed by comparing the number of skin tone macroblocks in a rectangular template to the number of skin tone macroblocks adjoining the top and sides of the template. There may also be provided spatial cross median filtering and segmentation to simplify the shape comparison.

For a better understanding of the present invention, together with other and further objects, reference is made to the following description, taken in conjunction with the accompanying drawings and its scope will be pointed out in the appended claims.

#### Brief Description of the Drawings

Figure 1 is a diagram illustrating steps of the preferred embodiment of the present invention.

Figure 2(a) is a plane of color space in the Cb-Cr chrominance coordinates of MPEG video.

Figure 2(b) is the figure 2(a) chrominance plane of color space illustrating the range of skin tones.

Figure 3(a) is a reproduction of a color image from video.

Figures 3(b) through 3(d) illustrate the skin color portions of the figure 3(a) image with varying skin tone threshold values.

Figure 4(a) is a reproduction of another color image.

Figure 4(b) illustrate the skin color portions of the Figure 4(a) image.

Figure 5 is a graph showing the effect of color classification threshold on false alarm rate and false dismissal rate.

Figure 6(a) depicts the identification of skin tone areas using NTSC chrominance signals of the Figure 3(a) color image according to the prior art for the Figure 3(a) color image.

Figure 6(b) depicts the identification of skin tone areas using MPEG chrominance signals according to a technique useful in the present invention.

Figure 7(a) depicts macroblock classification of the color image of Figure 4(a).

Figure 7(b) depicts the effect of spatial median filtering on the macroblock classification of Figure 7(a).

Figures 8(a), 8(b), 8(c) and 8(d) depict the application of template matching for shape comparison in accordance with the present invention.

Figure 9 is a diagram illustrating shape comparison in accordance with the present invention.

Figures 10(a), 10(b) and 10(c) illustrate the segmentation process useful in connection with the present invention.

Figure 11(a) is a reproduction of a color image.

Figure 11(b) is a segmented image identifying regions of skin tones for the Figure 11(a) image.

Figure 11(c) is an identification of face regions for the Figure 11(a) image following shape comparison.

Figure 12 illustrates groupings of DCT coefficients in accordance with the prior art.

Figure 13 illustrates inverse motion compensation.

Figure 14(a) and 14(b) illustrate the run time for a selection of frames using the method of the present invention on two workstations.

Figures 15(a) to (h) illustrate examples of face detection results for various video frames in accordance with the method of the present invention.

002250 4250550

Figures 16(a) and (b) illustrate the effects of threshold values in connection with the spatial frequency aspects of the present invention.

### Description of the Invention

The present invention provides a fast method that automatically detects  
5 face regions in MPEG-compressed video. In an exemplary embodiment, the  
method is applied to color video frames in the form of the inverse-quantized DCT  
coefficients of the MPEG macroblocks, and generates bounding rectangles of  
detected face regions. By detecting faces using just the DCT coefficients, we  
avoid the computational intensive inverse DCT transform. Thus, only minimal  
10 decoding of MPEG video is necessary, and the algorithm is able to achieve high  
speed.

MPEG video consists of three different types of frames, namely I  
(intra-frame coded), P(one-way predictive coded), and B(bi-directional predictive  
15 coded) frames. For the purposes of indexing and search, face detection in I frames  
is usually sufficient. This is because faces in video scenes usually stay much  
longer than the duration of an MPEG group of pictures (GOP), which usually  
consists of 12 to 15 frames (about 0.5 second duration). After minimal decoding  
of an MPEG stream, the DCT coefficients can be obtained easily for the luminance  
and chrominance blocks in I frames. If face detection in B or P frames is desired,  
20 one can apply transform-domain inverse motion compensation to obtain the  
corresponding DCT coefficients for blocks in B and P frames. The DCT  
coefficients of the translated block in the reference frame can be computed using  
the algorithm proposed by Chang and Messerschmitt. The idea is that the DCT  
coefficients of a translated and non-aligned block can be obtained by summing  
25 weighted DCT coefficients from their four overlapping neighbor blocks. The  
computation intensive and time consuming inverse DCT transform is not needed.  
The method becomes very efficient if only part of the DCT coefficients (e.g., DC  
coefficients) are used. By obtaining the DCT DC coefficients of P and B frames,  
the method can also be applied to these video image frames.

Figure 1 is a block diagram of the preferred embodiment of a face detection method of the invention. In the diagram, rounded rectangles represent input data, intermediate and final results; rectangles represent operations in the method. The method has three stages, where average chrominance of macroblocks, shape constraints on human faces, and energy distribution of the DCT coefficients are used respectively. MPEG macroblocks (16x16 pixels) are the preferred processing unit representing incremental image portions, so that the bounding rectangles of the detected face regions have a resolution limited by the size of the macroblocks. The result of the method is a list of face regions and their locations in the video image.

The method also uses domain knowledge to help make decisions at each stage. Domain knowledge is shown as ellipses in the diagram. Statistics of human skin-tone colors in the chrominance plane is used in Stage 1. We use sample patches of face regions and non-face regions as training data to generate the statistics. Shape constraints on human faces are applied in Stage 2. One of them is the anatomical constraint of human faces. For example, it is impossible for the outline of a human face to have an aspect ratio (height over width) of 3 to 1, or 1 to 3, if we do not consider face regions in video created by special effects. Other constraints are from the attributes of MPEG video. For example, the size of the video frames sets the upper bound of the largest face regions that our method can detect. In Stage 3, knowledge of the energy distribution over the DCT coefficients of face regions is used.

In the first step of the method, DCT DC values of Cb and Cr blocks are used to represent average chrominance of the corresponding macroblocks.

According to its average Cb and Cr values, each macroblock is classified as a skin-tone macroblock or not a skin-tone macroblock, based on the statistical distribution of skin-tone colors in the chrominance plane. After Stage 1, all macroblocks with skin-tone colors are considered as candidate face regions. Therefore, a binary mask image can be generated for each frame, in which a "one" means a candidate face macroblock, and a "zero" means the opposite. The binary

mask image is post-processed by morphological operations to eliminate noise and fill up holes in it.

It should be noted that the goal of the first step is to detect all candidate face blocks. A moderate level of false alarms is tolerable at this stage. Additional constraints in later stages can be used to greatly reduce the false alarm rate.

In the second step, our goal is to detect face regions in the mask images generated by Stage 1. We use a projection method to quickly find the target areas, then apply an iterative template matching procedure to locate the individual candidate face regions. To cope with various sizes and orientations of faces, shape constraints are used to eliminate false indications of face regions.

In the final stage, for each face region detected in Stage 2, we calculate the energy distribution of the luminance DCT coefficients over different frequency bands in the DCT domain. This is based on the observation that human faces contain uneven frequency components in different orientations. The result is used as a final verification of face detection and helps eliminate anomalies in Stage 2. Note that this stage can be omitted for P and B frames in order to save computation of DCT AC coefficients in these types of frames.

Overall, our algorithm has a cascading structure with various kinds of domain knowledge applied. To speed up the algorithm, our principle is to push simpler stages up to the beginning, and leave the most complex ones to the end. Thus, more complex stages only have to work on a subset of the original data so that computation is reduced.

In the first step of the method, we check each macroblock or other incremental portion of the video frame to see if it is a candidate face portion or not. The key to this classification is the uniqueness of human skin-tone colors. We use training data to generate skin-tone color statistics, then apply the Bayesian minimum risk decision rule to classify each image portion.

In video transmission and storage, colors are usually separated into luminance and chrominance components to exploit the fact that human eyes are less sensitive to chrominance variations. Psychophysical experiments indicate that

perception of colors has three attributes: hue, saturation, and intensity. Intensity corresponds to the luminance value (Y), while hue and saturation are kept in the chrominance components (such as Cr and Cb).

Human skin tones form a special category of colors, distinctive from the colors of most other natural objects. Although skin colors differ from person to person, and race to race, they are distributed over a very small area on the chrominance plane. This means that skin colors are relatively consistent in hue and saturation. The major difference between skin tones is intensity or luminance.

The above fact has been noticed and used by researchers in consumer electronics to design TV circuits that automatically detect and correct human skin-tone colors that are sensitive to human eyes. I and Q components of NTSC chrominance signals are used to estimate the hue and saturation of a color. Colors are classified as skin tones if their hue and saturation fall into certain ranges. By taking out the luminance component of colors, the difference between skin colors of different races and the effect of lighting conditions are reduced.

In our work, we make use of human skin-tone characteristics as well, but take a different approach. In connection with MPEG video, we use Cb and Cr, instead of I and Q, because the former are usually the chrominance components used in MPEG video. Second, we generate statistics of skin-tone color distribution directly on the Cb-Cr plane and use it in the classification process. Thus, we can get more accurate statistics and avoid the non-linear transform to get hue and saturation. Third, to enhance the robustness and flexibility of our algorithm, we classify the colors based on the Bayesian decision rules.

Figure 2(a) shows the distribution of all displayable colors in the RGB color cube on the Cr-Cb chrominance plane. Conversion between R, G, B and Y, Cr, Cb is as follows:

$$\begin{bmatrix} Y \\ Cb \\ Cr \end{bmatrix} = \begin{bmatrix} 0.299 & 0.587 & 0.114 \\ -0.169 & -0.331 & 0.500 \\ 0.500 & -0.419 & -0.081 \end{bmatrix} \cdot \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (1)$$

Note that the colors in Figure 2(a) are for illustration purposes only, because intensity information cannot be revealed on this chrominance graph. One color in this figure actually corresponds to all the colors that have the same chrominance values but different intensities.

- 5 To generate skin-tone color statistics on the Cb-Cr plane, we use as training data sample face patches of different races, which are cut from video frames. Figure 2(b) shows the distribution of skin-tone colors in the Cr-Cb plane, based on over forty sample face patches. The white area in Figure 2(b) corresponds to chrominance values that have a non-zero probability in the sample face patches.
- 10 Compared with Figure 2(a), it is clear that skin-tone colors are distributed over a very small area in the chrominance plane.

- We use the Bayesian decision rule for minimum cost to classify a color into skin-tone class or non-skin-tone class. This technique is flexible because it allows us to use the statistics of skin-tone colors, and to take into consideration the
- 15 different effects of false alarms and false dismissals. The Bayesian decision rule for minimum cost (for two classes) is described as follows.

$$R_0(X) = C_{00} \cdot p(\omega_0|X) + C_{10} \cdot p(\omega_1|X) \quad (2)$$

$$R_1(X) = C_{01} \cdot p(\omega_0|X) + C_{11} \cdot p(\omega_1|X) \quad (3)$$

$$R_0(X) < R_1(X) \Rightarrow X \in \omega_0 \quad (4)$$

$$R_0(X) > R_1(X) \Rightarrow X \in \omega_1 \quad (5)$$

- $\omega_0$  and  $\omega_1$  denote two classes, respectively.  $p(\omega_i|X)$  denotes the *a posteriori* probability, i.e., the probability of being in class  $i$  given sample  $X$ .  $C_{00}$  and  $C_{11}$  denote the cost coefficients of correct classifications;  $C_{01}$  and  $C_{10}$
- 20 denote the cost coefficients of false classifications. Therefore,  $R_i(X)$  is the "cost" of classifying an unknown sample into class  $i$ . The classification problem becomes finding the class which gives the minimal cost, considering different cost weightings on classification decisions.

In our context, the two classes are non-skin-tone class ( $\omega_0$ ) and skin-tone class ( $\omega_1$ ). We assign zero cost to correct classifications, so that  $C_{00}$  and  $C_{11}$  are both zero. Then the minimum cost decision rule reduces to the following.

$$C_{10}p(\omega_1|X) < C_{01}p(\omega_0|X) \Rightarrow X \in \omega_0 \quad (6)$$

$$C_{10}p(\omega_1|X) > C_{01}p(\omega_0|X) \Rightarrow X \in \omega_1 \quad (7)$$

Applying the Bayesian Formula (Eq. (8)) to the above equations, we obtain the decision rules used in our algorithm as shown by Eqs. (9) and (10).

$$p(\omega_i|X) = \frac{p(X|\omega_i)p(\omega_i)}{p(X)} \quad (8)$$

$$\frac{p(X|\omega_1)}{p(X|\omega_0)} < TH \Rightarrow X \in \omega_0 \quad (9)$$

$$\frac{p(X|\omega_1)}{p(X|\omega_0)} > TH \Rightarrow X \in \omega_1 \quad (10)$$

$$\text{where } TH = \frac{C_{01}}{C_{10}} \cdot \frac{p(\omega_0)}{p(\omega_1)}$$

In the above equations,  $p(\omega_i)$  is the corresponding *a priori* probability of class  $\omega_i$ .  $p(X|\omega_i)$  denotes the conditional probability density functions of skin or non-skin colors on the chrominance (Cb-Cr) plane. The conditional probabilities are generated by the method described above using sample face and non-face patches as training data.

TH is the adjustable decision threshold. The higher  $C_{10}$  (or the lower  $C_{01}$ ) is, the more false alarms are allowed, and vice versa. By changing  $C_{10}$  and/or  $C_{01}$ , we can control the amount of false alarms and false dismissals allowed in Stage 1. In applications, TH set at 2.0 is a reasonable value, based on our experiments of the effect of TH on the classification.

Figure 3 demonstrates the effect of the classification threshold (TH) on the classification of color pixels. A video frame with complex background and multiple small faces is shown in Figure 3(a). As TH decreases, more skin-color pixels are detected, which are shown as gray pixels in Figures 3(b), 3(c), and 3(d).

For a complex scene like Figure 3(a), it is clear that a small TH value is needed to generate solid clusters of face pixels for the detection of face regions.

Figure 4(a) shows a video frame with a relatively clean background and a large human face. The result of pixel classification with TH equals 0.5 is shown in Figure 4(b). When TH is such a small value, we can detect almost all of the face pixels, with some noise in the background. Figure 5 illustrates the variations of false detection rate and false alarm rate when TH changes from 1 to 20. We use the classification result of Figure 4(b) as the basis for comparison. First, we manually segment Figure 4(a) into the face region (a bounding rectangle) and the background region. Then, in the result mask images, for each TH value, we count the number of "one" pixels in the face region, and the number of "one" pixels in the background. These numbers are compared with those for TH equals 0.5, to derive the relative false alarm rate, and relative false dismissal rate. Figure 5 shows that as we raise the classification threshold, the false alarm rate decreases.

However, at the same time, we have more false dismissals. Because false dismissals are undesirable at the early stage of face detection, a TH value of 2.0 is reasonable for classification. Nevertheless, it can be changed in order to have a higher detection rate or a lower false alarm rate in applications.

Experiments show that our method based on the statistical decision rule gives more accurate results, compared with the method suggested by Rzeszewski. An example of comparison is given in Figure 6. Figures 6(a) and 6(b) are the classification results using the prior art method and our method, respectively. Gray pixels correspond to skin-tone colors; black pixels the opposite. The original video frame is Figure 3(a). It can be seen that Figure 6(b) has more solid clusters of candidate face regions.

We apply the above minimum cost decision rule to MPEG video streams, and classify each MPEG macroblock or other image portion as a candidate face macroblock or a non-face one. We use only the DCT DC coefficients of the corresponding Cr and Cb blocks, which are equivalent to the average values (up to a scale) of the chrominance blocks in the pixel domain. Higher resolution in

skin-tone detection can be achieved by taking more DCT coefficients (e.g., a few low-order coefficients) as input in the above classification process.

After the classification, we get a binary mask image for each video frame. Each value in the mask indicates the classification results of the corresponding macroblock. Then, a 3x3 (macroblocks) cross median filter is applied to the binary mask image to remove noise and smooth the image. The filtering helps because faces are connected regions, and are homogenous in chrominance. We use the video frame of Figure 4(a) as an example. Figure 7(a) shows the binary mask image after the classification by average chrominance. Figure 7(b) is the better result after applying the median filter to Figure 7(a).

As shown in Figure 1, after Stage 1 of the method, we have a macroblock mask image for each video frame. In these macroblock mask images, a "one" pixel corresponds to a macroblock whose average color is a skin-tone color. Now our task is to scan through these mask images, and detect actual face regions in them.

Clearly, chrominance information alone is not enough to detect face regions. In a video sequence with complex scenes, besides human faces, there may be other exposed parts of the body with skin tones, and natural scenes with colors similar to skin tones (e.g., a desert scene). All these examples would produce positive yet false detections in Stage 1 of our method. In Stage 2, we apply shape constraints of human faces on the binary mask images generated by Stage 1, to eliminate these false alarms and detect candidate face regions.

Just like color, the shape of human faces is unique and consistent. The outline of a human face can be approximated by an ellipse, or more precisely, by connected arcs. Furthermore, typical face outlines have been found to have aspect ratios in a narrow range between 1.4 and 1.6, and tilt in the range (-30, +30) degrees.

Our face detection method works on the macroblock level, so that it is difficult to use ellipses or arcs to describe face outlines. As mentioned previously, this is a trade-off for high speed. As an approximation, rectangles with certain aspect ratios can be used as the boundary of face regions. We set the range of

aspect ratios of these bounding rectangles to  $[1, 1.7]$ . It is a larger range, because we are using rectangles to estimate the aspect ratio of the actual face outline.

Besides certain aspect ratios, these rectangles are also bounded by size.

The size of the video frames upper bounds the size of face regions. It is lower

5 bounded as well, because it is generally believed in the face recognition field that  $32 \times 32$  pixels is the lower limit for face detection. Since we are working in the compressed domain, we set the lower limit of our face detection method to  $48 \times 48$  pixels, or,  $3 \times 3$  macroblocks. Faces smaller than this size are not considered for detection.

10 In summary, the shape constraints we put on macroblock mask images are as follows: (1) faces are contiguous regions that fit well in their bounding rectangles, whether the face is front view or side view, or whether the head is upright or a little tilted; (2) the size of the bounding rectangles is lower bounded by the lower limit of face detection, and upper bounded by the size of the video  
15 frames; (3) the aspect ratios of the bounding rectangles should be in a certain range.

Stage 1 of the method generates candidate face macroblocks after skin-tone classification. Compared with original video frames, the resolution of the mask images is 16 times lower in horizontal and vertical directions. Existing geometric  
20 analysis techniques, e.g., those detecting arcs corresponding to chin boundaries and hairlines, are not appropriate at this resolution. Note that we may detect all candidate regions back in the pixel domain and apply the above method. But to keep our method in the compressed domain, we take into account the resolution limit and modify the task objective. Now our task is to find contiguous regions in  
25 the macroblock mask images that can be bounded by a rectangle satisfying some size and shape constraints.

To accomplish this task, we use a method called binary template matching. The idea is to use rectangles of possible sizes and aspect ratios as face templates to match against the binary mask images. A face template is shown in Figure 8(a),  
30 whose size is  $(M+1) \times (N+1)$  macroblocks. Note that internally, the face templates

are represented by only  $(M+1) \times (N+1)$  pixels. For the sake of illustration, mask images are blown up in the figures by a factor of 16 in each direction. The template consists of two parts: the face region, which is the  $M \times N$  shaded rectangle, and the background, which is the area between the inner and outer rectangles. The size of the face region ( $M \times N$ ) should be bounded by size and aspect ratio, according to our shape constraints. The reason we consider the background as part of the template is that the color of the background adjacent to the face region is usually distinctive from skin tone, so that there should be few "ones" in this region. The macroblocks adjacent to the bottom of the face region are not considered in the template, because they can be either the exposed neck or clothes and have no definitive color characteristics.

The matching criterion is two-fold. As we slide this two-frame template over a macroblock mask image, we count both the number of ones covered in the shaded region, and the number of ones in the background region. The intuition is that for a match, the first number should be high, and the second should be low. Since we are processing binary mask images, no multiplication is involved in the matching. Therefore this procedure is of low complexity. We count the number of ones inside the face (shaded) rectangle, as well as the numbers of ones in the top, left, and right parts of the background region. Denote these numbers as  $N_0$ ,  $N_1$ ,  $N_2$ , and  $N_3$ , respectively. Only if  $N_0$  is above a threshold, and  $N_1$ ,  $N_2$ ,  $N_3$  are below certain thresholds, do we declare a match.

This is illustrated in Figure 8. Figure 8(b) is a match, because the face region is almost covered by ones, and there are few ones in the background region. Figure 8(c) is not a match, because the face rectangle is not covered enough by ones. Figure 8(d) is not a match either, because there are too many ones in the background region.

The flow chart of Stage 2 of our algorithm is shown in Figure 9. The core of this stage is the binary template matching technique described above. To further speed up the processing of this stage, extra work is done to limit the search area for template matching.

For each macroblock mask image, we first segment it into non-overlapping rectangular regions that contain either no ones or a contiguous one region. This segmentation is done by projecting the mask image onto the x and y axes. Given a binary mask image  $B(x,y)$  ( $x=0,1,\dots,M-1$ ;  $y=0,1,\dots,N-1$ ), the projections are defined as follows, where PX denotes the projection onto the X axis; PY denotes the projection onto the Y axis.

$$PX(x) = \sum_{y=0}^{N-1} B(x,y) \quad x=0, 1, \dots, M-1 \quad (11)$$

$$PY(y) = \sum_{x=0}^{M-1} B(x,y) \quad y=0, 1, \dots, N-1 \quad (12)$$

Based on the zero-runs and non-zero-runs in PX and PY, we are able to segment the binary mask image. As an example, Figure 10(a) is the binary mask image corresponding to a video frame with two faces in it. After the projections, we have  $PX=[0, 0, 0, 0, 0, 0, 5, 6, 4, 2, 0, 0, 0, 2, 3, 3]$ , and  $PY=[0, 0, 0, 0, 4, 6, 6, 4, 3, 2, 0, 0, 0, 0, 0]$ . Using the above projection results, the binary mask image is segmented as shown in Figure 10(b). The way we segment the image guarantees that in each of the segments, there is either a contiguous one region, or no ones at all.

For each of the segments, we first count the number of ones in the segment. If the number is zero, or less than the minimum value for a possible face region, there is no need to perform binary template matching in it, so we simply proceed to the next segment, as shown in the flow chart in Figure 9. Therefore, for all the rectangular segments shown in Figure 10(b), only the two regions in Figure 10(c) need to be searched for faces.

Then, in each of these segments, we apply the binary template matching method. Because the size of the face region is unknown, we start from the largest possible rectangle for each segment, then gradually reduce the template size. Therefore, all sizes of the face regions can be detected.

Finally, overlapping face bounding rectangles are resolved. If only a small area is overlapped, this may be a situation where two faces are very close to each other. Therefore, we keep both regions as valid face regions. If one of the regions is small and the overlapping area is large compared with its size, we discard the smaller rectangle.

An example is shown in Figure 11. The original video frame is in Figure 11(a). The binary mask image is in Figure 11(b), along with search regions (bounded by white rectangular frames). Figure 11(c) shows the detected face regions before Stage 3. The final result after Stage 3 is overlaid on Figure 11(a). In some cases, there might be non-face regions that have similar colors to skin tones, and have a roughly rectangular shape. This will cause false alarms, such as the rectangle in the lower-left corner of Figure 11(c). This problem can be solved in Stage 3 of our algorithm, as will be described below.

The main purpose of the last stage of our method is to verify the face detection result generated by the first two stages, and remove false alarms caused by objects with colors similar to skin tones. Because of the existence of eyes, nose-mouth junction, and lips in face regions, there are many discontinuities of intensity level in the vertical direction of the image in face regions. These discontinuities correspond to the DCT coefficients in the high vertical frequency area. Therefore, we expect some energy in the luminance DCT coefficients (Y-component in MPEG video) in that frequency band. This is the rule we use in this stage for verification. Before calculating the energy, we have to group the 64 DCT coefficients into different frequency bands.

Various methods have been proposed in the signal processing field to classify DCT coefficients into groups of different spatial frequencies. In our algorithm, we follow the DCT coefficient grouping scheme proposed by Ho and Gersho for vector quantization of transformed images, where the DCT coefficients in a 8x8 transform block are partitioned into groups corresponding to the directional features (horizontal, vertical, and diagonal edges) of the block in the spatial domain, along with the DC coefficient. Figure 12 shows the grouping of

the DCT coefficients, where groups H, V, and D correspond to vertical, horizontal, and diagonal edges, respectively. In these matrices, one means the DCT coefficient at that position belongs to the group, and zero means the opposite. In our algorithm, we use only the H and V matrices, along with the DCT DC coefficients of luminance blocks.

Given a candidate face region of size  $M \times N$  macroblocks, we compute the energy of the corresponding luminance blocks in the DC and H, V areas as follows:

$$E = \sum_{i=0}^{M \times N \times 4 - 1} \left( \sum_{m=0}^7 \sum_{n=0}^7 |DCT_i(m, n)|^2 \right) \quad (13)$$

$$E_{DC} = \left( \sum_{i=0}^{M \times N \times 4 - 1} |DCT_i(0, 0)|^2 \right) / E \quad (14)$$

$$E_H = \left( \sum_{i=0}^{M \times N \times 4 - 1} \left( \sum_{m=0}^7 \sum_{n=0}^7 |DCT_i(m, n) \cdot H(m, n)|^2 \right) \right) / E \quad (15)$$

$$E_V = \left( \sum_{i=0}^{M \times N \times 4 - 1} \left( \sum_{m=0}^7 \sum_{n=0}^7 |DCT_i(m, n) \cdot V(m, n)|^2 \right) \right) / E \quad (16)$$

Eq. (13) shows that  $E$  is the summation of the energy of all the DCT coefficients in the candidate face region. It equals the energy of the pixel values of this face region, because of the DCT transform's energy-conserving property.  $E_{DC}$ ,  $E_H$ ,  $E_V$  are the normalized energies of all the DCT coefficients in the candidate region, of groups DC, H, and V, respectively. Note that we assume 4:2:0 macroblock structure for the MPEG video, so that in a region of  $M \times N$  macroblocks, there are  $M \times N \times 4$  DCT luminance blocks. Matrices H and V are given as in Figure 12.

We set up two thresholds  $T_{DC}$  and  $T_{V/H}$ . If either  $E_{DC} > T_{DC}$ , or  $E_V/E_H < T_{V/H}$ , we declare the face region is a false alarm. The reason is that a face region should not have near 100% energy in DC value, because face regions contain details and edges. Also the energy corresponding to the horizontal edges ( $E_V$ ) should be large enough. We use  $E_V/E_H$  instead of absolute  $E_V$ , because face regions may have different resolutions in an unconstrained video stream, so that  $E_V$  may have a large

dynamic range. Therefore, the ratio of  $E_V/E_H$  is more consistent and reliable than the absolute value of  $E_V$ .

Using these thresholds, we verify each candidate face region declared by Stage 2 of our algorithm. This helps remove some false alarms from Stage 2. For example, the false detection in Figure 11(c) is removed after the verification, and does not appear in the final detection result (Figure 11(a)). The effect of the  $T_{DC}$  and  $T_{V,H}$  on the performance of our algorithm will be discussed below.

In MPEG video streams, I frames are intra-coded. After minimal parsing of an MPEG stream, all the DCT coefficients of an I frame are available for its luminance and chrominance blocks. Thus, our algorithm can be applied directly to MPEG I frames, using the DCT DCs of chrominance blocks and the DCT coefficients of luminance blocks.

P frames consist of motion compensated (MC) macroblocks and intra-coded macroblocks. For an intra-coded macroblock, all its DCT DCs can be obtained directly. An MC macroblock is coded using a motion vector and the DCT transformed residue errors. Each macroblock consists of four luminance blocks and two chrominance blocks (for 4:2:0 chrominance format), namely Cb and Cr blocks. Using partial inverse motion compensation, the DCT DC values of P frames can be computed efficiently, without the inverse DCT transform. Here we use the efficient method proposed by Meng et al.

To apply our algorithm, we need to compute the DCT DC coefficients of Cb, Cr blocks in P frames. In practical videos, variance within each chrominance block is small. Thus, the DCT DC of a MC block in P frame can be approximated by taking the area weighted average of the four blocks in the previous reference frame pointed to by the motion vector. This is shown by Eq. (17) and Figure 13, where  $x$  and  $y$  are the horizontal and vertical components of the motion vector, modulo block size 8;  $b_0, b_1, b_2$ , and  $b_3$  are the DCT DC coefficients of the four neighboring blocks pointed to by the motion vector;  $b_{error}$  is the DCT DC value of the MC residue error of the block to be computed;  $b$  is the inverse motion

compensated DCT DC value. The layout of the blocks and the motion vector are illustrated in Figure 13.

$$b = [b_0*(8-x)*(8-y) + b_1*x*(8-y) + b_2*(8-x)*y + b_3*x*y] / 64 + b_{error} \quad (17)$$

Note that motion estimation is usually done on luminance macroblocks. The motion vectors for the chrominance blocks has to be adjusted according to the encoding chrominance format. For example, in a 4:2:0 chrominance scheme, the motion vector has to be reduced by half before using Eq. (17).

Inverse motion compensation can be used to reconstruct DCT AC coefficients in P frames as well. However, this is more expensive than just computing the DCT DCs. Therefore, for P frames, we use only the first two stages of our algorithm to detect face regions. This will not affect the detection rate. However, more false alarms are expected because of the omission of the last verification stage.

B frames in MPEG streams are similar to P frames, except that they are bi-directionally motion compensated. For each MC macroblock, either forward MC, or backward MC, or both are used. The technique above can still be applied.

We used 100 I-frames from a 10-minute MPEG-compressed CNN news video as our first test set. These frames include anchorperson scenes, news stories, interviews, and commercials which cover video scenes with various complexities. The number of each category of frames is shown in Table 1 to briefly summarize the content of this test set. The size of each frame is 352x240 pixels. 4:2:0 macroblock format is used, which means that chrominance signals are subsampled by two, both horizontally and vertically.

002260-4290560

5

Category	Number
Anchor	15
Commercial	10
Interview	15
New Story	60
Total	100

There are 91 faces in these frames, including frontal views, semi-frontal views, side views, and tilted faces. These faces are often in complex scenes. Some of the frames have multiple faces. A detailed description of the faces is shown in Table 2.

10

15

Category	Number
Frontal	61
Semi-Frontal	13
Side	11
Tilted	6
In Multiple Faces	37
Number of Persons	15
Min Face Size	50x50
Max Face Size	240x240

20

We also use 50 P frames from the same news video as test set 2, to test our face detection method in inter-coded frames. There are 44 faces in these P frames.

25

Finally, we use another 100 I frames from other CNN news clips as test set 3. It includes anchorperson scenes (with different anchors from those in test set 1), two news stories, and one out-door interview. There are 46 faces in these frames.

The run time of our algorithm varies depending on the content of the video frames. Stage 1 of the algorithm treats all incoming video frames equally, and is run on each macroblock. So run time for Stage 1 is fixed. The difference is in Stages 2 and 3.

5       The run time of Stage 2 depends on the output of Stage 1. If after Stage 1, very few ones are detected in the macroblock mask image, then little binary template matching is involved, so less time is needed. On the contrary, if the video scene contains many skin-tone color regions, and is complex, Stage 2 will spend a longer time.

10       The run time of Stage 3 is proportional to the area of face regions detected in Stage 2, because, for each macroblock, we have to calculate its energy distribution in the DCT domain.

Overall, based on experiments, the speed of our algorithm is as follows. On a SPARC 5 workstation, the average run time is 32.6 milliseconds for the  
15   100 I frames in test set 1. On a SGI Indigo 2 workstation, the average run time is 15.6 milliseconds. We measure the run time of the algorithm using C library functions of timing with the precision of microseconds. No inverse quantization is involved. Figure 14 shows the histogram of run time of our algorithm on the 100 frame test set 1. The run times on both SPARC and SGI workstations are  
20   given in Figures 14(a) and 14(b), respectively.

For P frames, since we only carry out the first two stages of the algorithm, the run time is shorter. Based on the experiments on the 50 P frames in test set 2, the average run time is 13.4 milliseconds on a SPARC 5 workstation; 7.0 milliseconds on a SGI Indigo 2 workstation.

25       For test set 1, our method detects 84 of the faces (92%), including faces of different sizes, frontal and side-view faces, etc. Detected face regions are marked by white rectangular frames overlaid on the original video frames. There are 8 false alarms in our experiment on test set 1.

For test set 2 (P frames), the algorithm detects 39 of the faces (88%), with  
30   15 false alarms. The false alarm rate is higher because Stage 3 of the algorithm is

skipped for inter-coded frames. However, as most face regions last longer than the duration of a GOP, and continue to appear in successive frames, analysis of continuation over time can be used to improve the accuracy of face detection in P and B frames. For test set 3, the algorithm detects 39 of the faces (85%), with 6  
5 false alarms.

Figure 15 shows some examples of the result of our face detection algorithm, from all of the test sets. These examples include frames with one or more faces of different sizes, frontal and side-view faces, and a frame without a face (Figure 15(g)). No face region is detected for Figure 15(g), although there are  
10 exposed arm and hand areas that have skin-tone colors.

As mentioned above, the thresholds used in Stage 3 affect the performance of the algorithm. In our experiment, we set  $T_{DC}=90\%$ ,  $T_{V/H}=0.5$ . Without Stage 3, our detection rate would be the same, but the number of false alarms would increase to 18 in the case of test set 1. Figures 16(a) and (b) respectively show the  
15 effect of  $T_{DC}$  and  $T_{V/H}$  on the number of false alarms and false dismissals in test set 1, where the thresholds are applied separately. In the real algorithm, the two thresholds are combined to achieve a better performance.

Figure 16 shows that, as we raise the  $T_{DC}$ , the number of false alarms will increase, and the number of false dismissals will decrease. If we raise  $T_{V/H}$ , the  
20 opposite will happen. Note that in Figure 16, the number of false alarms corresponds to that of the algorithm, while the number of false dismissals corresponds to that in Stage 3 only. The false dismissals in the first two stages of the algorithm cannot be recovered by Stage 3. From Figure 16 we also see that the thresholds we use for the algorithm are reasonable.

25 The present method can only be applied to color images and videos, provided chrominance information is available for use in Stage 1 of the method. When macroblocks are used, the smallest faces that are detectable by this algorithm are about 48 by 48 pixels (3 by 3 macroblocks), bounded by the lower limit of machine detection of faces, and the fact that for MPEG, the method can

operate in the compressed domain. The method is relatively independent of lighting conditions, but very poor lighting conditions still cause false dismissals.

False dismissals cannot be totally avoided, especially in very cluttered scenes with many small faces (e.g., a scene in a football game). This situation would be alleviated if we use video sequences with larger frame sizes. False alarms usually happen because of the existence of regions that have skin-tone colors, but which are not human faces, for example: desert/soil scene, yellowish and reddish light, etc. There are fewer false alarms after we apply the shape and energy constraints. By adjusting the thresholds in Stages 1 and 3 of our method, we can make face detection more or less conservative, i.e., to have higher detection rate with higher false alarm rate, or the opposite.

The method of the invention is efficient and can be applied to large video databases for indexing and recognition. It helps focus our attention to only a small portion of the entire video sequences and frames. Once we detect these target regions, we can decode them back to the pixel domain, in which more sophisticated techniques can be applied, to either enhance or verify our face detection results, or apply video indexing or face recognition techniques.

The method does not give the exact outlines of the faces, when we avoid inverse DCT transform and work at the macroblock resolution. The positions of the faces detected are sometimes not perfectly aligned, because the face rectangles we detect lie on the borders of 16x16 macroblocks. This can be improved if the compressed video sequence has a format with more chrominance information, e.g., the 4:2:2 format, so that we can work on 8x8 blocks and the face detection result will be improved. Also, we can use more DCT coefficients rather than just the DC coefficient of the Cb and Cr blocks, to get more accurate results.

While there has been described what is believed to be the preferred embodiment of the invention, those skilled in the art will recognize that other and further changes may be made thereto without departing from the spirit of the invention, and it is intended to claim all such changes and modifications as fall within the true scope of the invention.